

# Using Meta-Knowledge Mined from Identifiers to Improve Intent Recognition in Conversational Systems

Claudio S. Pinhanez, Paulo Cavalin, Victor Ribeiro, Ana Paula Appel, Heloisa Candello, Julio Nogima, Mauro Pichiliani, Melina Guerra, Maira Gatti de Bayser, Gabriel Malfatti, Henrique Ferreira

IBM Research - Brazil  
csantosp@br.ibm.com

## Abstract

In this paper we explore the improvement of intent recognition in conversational systems by the use of meta-knowledge embedded in intent identifiers. Developers often include such knowledge, structure as taxonomies, in the documentation of chatbots. By using neuro-symbolic algorithms to incorporate those taxonomies into embeddings of the output space, we were able to improve accuracy in intent recognition. In datasets with intents and example utterances from 200 professional chatbots, we saw decreases in the equal error rate (EER) in more than 40% of the chatbots in comparison to the baseline of the same algorithm without the meta-knowledge. The meta-knowledge proved also to be effective in detecting out-of-scope utterances, improving the false acceptance rate (FAR) in two thirds of the chatbots, with decreases of 0.05 or more in FAR in almost 40% of the chatbots. When considering only the well-developed workspaces with a high level use of taxonomies, FAR decreased more than 0.05 in 77% of them, and more than 0.1 in 39% of the chatbots.

## 1 Introduction

Classification of sentences into a discrete set of classes is a key part of professional conversational systems. In fact, most of those systems require developers to define the different classes, or *intents*, by enumerating exemplars of each of them, since classification is often performed using *machine learning (ML)* methods. The process of classifying an input sentence into a specific intent or signaling it as *out-of-scope (OOS)* of the system is often referred to as *intent recognition*.

Determining a class solely on a list of exemplars is a practical method to implement ML systems but it is hardly a natural way for human beings to define a class. In real life, people define classes often using a rich mix of symbolic definitions, sometimes

taxonomic in nature, such as in “a credit card is a type of bank card”, coupled with its sub-classes, for instance, “basic”, “premium”, and typical features such as “international”. People also use exemplars, “card X of bank Y is a credit card”, as well as particular examples to describe a sub-class, such as in “card W is an international card”. They also use counter examples, either categorically, “a debit card is not a credit card”, or in examples, “card Z is not a credit card”. Defining and specifying classes in the real world is, in fact, a cultural, contextual, and linguistic construct, and how people and societies perform this process is a traditional research subject in social sciences, notably in anthropology (Durkheim and Mauss, 1963; Needham, 1979; Bowker and Star, 2000).

This paper explores algorithms for intent recognition which use both the sets of exemplars and taxonomic-like symbolic descriptions of a class to define and train intents in conversational systems using ML methods. We aim not only to provide methods more aligned to everyday class definition practices of developers but also to improve the accuracy of the ML methods. Inspired by a reverse dictionary algorithm (Kartsaklis et al., 2018) and previous work on keyword-based classification (Cavalin et al., 2020), we propose three *neuro-symbolic* algorithms which combine taxonomic descriptions of classes with traditional exemplar-based supervised learning. We show that those novel algorithms are able to decrease error rates for a significant number of datasets, particularly in the difficult task of detecting OOS cases in real, professional chatbots.

The key idea behind our algorithms is to substitute the typical *softmax* used in the output layer of a ML text classifier with a space of embeddings of the taxonomic descriptions of the intents. The training process uses the exemplars in standard ways while the recognition process is performed using similarity distances in the embedding output space.

This is similar to ideas used in *zero-shot learning* methods (Palatucci et al., 2009; Socher et al., 2013; Akata et al., 2015, 2016), in which classes defined by sub-concepts are also encoded with special embeddings to allow detection of new classes without exemplars.

We tested our algorithms using real datasets by exploring a common practice among developers of conversational systems, who often embed symbolic knowledge as documentation in intent identifiers. In a previous work (Pinhanez et al., 2021), we observed a pattern among developers of using taxonomic-like structures to name the intents in which strings of reoccurring concepts are used to identify and document the different classes. For example, an intent about utterances where users ask for the balance of a credit card may be named “checking\_credit\_card\_balance”, while an intent related to finding out the date of payment of the balance could be identified as “asking\_credit\_card\_balance\_payment\_date”.

We call those structures *intent proto-taxonomies*, and real examples are shown in figures 1 and 2. In (Pinhanez et al., 2021), we studied the use by developers of intent proto-taxonomies quantitatively and qualitatively, as well as proposed an algorithm to mine this meta-knowledge automatically, and concluded that their use is fairly common in at least one professional chatbot development platform. This paper focuses on the algorithms to use the meta-knowledge and on evaluating their impact on the accuracy of intent recognition.

The paper starts by looking into the recent advances in neuro-symbolic systems and describing briefly the practice of developers of conversational systems of embedding meta-knowledge within the source code of their systems. We follow by describing the proposed three algorithms integrating such meta-knowledge into intent recognition ML algorithms and by evaluating them first with two typical intent recognition datasets, and then with hundreds of workspaces created in a professional tool called here *ChatWorks*<sup>1</sup>. The results show most of those workspaces can benefit from the techniques described in this paper, notably for OOS detection tasks, often with accuracy improvements of 5% or more solely derived from the use of the additional symbolic description from the documentation.

---

<sup>1</sup>We use an anonymous name for the tool due to publication restrictions from the platform company.

## 2 Related Work

The value and limits of symbolic categorization in AI have been of interest since the early days (Newell, 1973; Richards, 1982; Kosslyn, 2006). But our work fits more in the context of a growing belief that symbolic knowledge needs to be included in ML systems, materialized in the so called *neuro-symbolic* approaches (Parisotto et al., 2017; Besold et al., 2017; Tenenbaum et al., 2011; Bengio, 2017; Mao et al., 2019; Hudson and Manning, 2019a; De Raedt et al., 2019).

Neuro-symbolic methods “*aim to transfer principles and mechanisms between (often nonclassical) logic-based computation and neural computation*” (Besold et al., 2017). Such kind of systems are viewed by some researchers as a way to embed high-level knowledge and even some form of “consciousness” into machine learning systems, making the language to develop them closer to “*what passes in a man’s own mind*” (Bengio, 2017).

In recent years, AI has witnessed a myriad of novel neuro-symbolic techniques and their application to different problems, contexts, and scenarios (Parisotto et al., 2017; Manhaeve et al., 2018; d’Avila Garcez et al., 2019; Hudson and Manning, 2019b; De Raedt et al., 2019). For instance, in (Mao et al., 2019), an approach for image understanding is suggested which takes the object-based scene representations and translates sentences into executable, symbolic programs. In (Oltamari et al., 2020), embeddings of knowledge graphs are used as attention layers for tasks such as autonomous driving (AV) and question-answering. And in (Kartsaklis et al., 2018), random walks in a knowledge graph are mapped as sentence embeddings for use in an *inverse dictionary problem*.

One important requirement for many neuro-symbolic systems is to represent knowledge in a structured format such as knowledge graphs, ontologies, or taxonomies (Ji et al., 2020). In some cases, such as the scene ontology for autonomous vehicles in (Oltamari et al., 2020), a lot of effort was needed for manual annotation. Nevertheless, as presented in (Fossati et al., 2015), an unsupervised approach can sometimes be used to mine the meta-knowledge introduced by the experts, such as the classes in *Wikipedia* pages.

Considering our context, intent identifiers are sometimes described using high-level representations of the class as we detail later. This is similar to what is used in some *zero-shot learning*



in (Pinhanez et al., 2021), it is possible to compute the *taxonomy rate* of a workspace by calculating the ratio between the number of intents with taxonomies and the total number of intents.

In 3,840 professional workspaces in the English language, it was found that 76% of them had a taxonomy rate above 10%, almost 52% had a taxonomy rate above 50%, and 16% had a very high taxonomy rate, above 90%. Moreover, the distribution followed a sort of “step” function where, as the threshold of 32 in the number of intents in a workspace was crossed, the majority of the workspaces had a taxonomy rate of more than 50%. It seems that, as the complexity of the workspace increases with the number of intents, more often developers resort to document them using an intent proto-taxonomy (see appendix B for details).

The use in our work of the intent proto-taxonomies as a symbolic description of classes is feasible because: (1) they are part of the documentation of the conversational system, so there is no need of acquiring knowledge from experts; (2) they are easily mined, as described in appendix A.

## 4 Using Taxonomic Intent Descriptions to Improve Intent Recognition

We present now a formal description of the methodology employed in this work which takes advantage of the intent proto-taxonomies using a neuro-symbolic approach. It expands some previous work which focused on the use of keywords as the source of symbolic information (Cavalin et al., 2020).

### 4.1 Embedding the Set of Classes

An *intent classification* method is a function  $D$  which maps a set of sentences (potentially infinite)  $S = \{s_1, s_2, \dots\}$  into a finite set of classes  $\Omega = \{\omega_1, \omega_2, \dots, \omega_c\}$ :

$$D : S \rightarrow \Omega \quad D(s) = \omega_i \quad (1)$$

To enable a numeric, easier handling of the input text, an embedding  $\xi : S \rightarrow \mathbb{R}^n$  is often used, mapping the space of sentences  $S$  into a vector space  $\mathbb{R}^n$ , and defining a classification function  $E : \mathbb{R}^n \rightarrow \Omega$  such that  $D(s) = E(\xi(s))$ . In most intent classifiers,  $E$  is composed of a function  $M$  which computes the likelihood of  $s$  being in a given class, often a neural network, followed by some sort of *argmax* function. Typically, *softmax* + *argmax* is used, noted simply as *softmax* here:

$$S \xrightarrow{\xi} \mathbb{R}^n \xrightarrow{M} \mathbb{R}^c \xrightarrow{\text{softmax}} \Omega \quad (2)$$

This paper explores how to use embeddings in the output side of the classification function, that is, by embedding the set  $\Omega$  of classes into another vector space  $\mathbb{R}^m$ , in some ways resembling the combination of object-based recognition and symbolic programming in (Mao et al., 2019). Instead, we combine here standard intent recognition methods with an encoding of taxonomies in knowledge graph-like structures. The idea is to use class embedding functions which somehow capture the knowledge in the intent proto-taxonomies.

Formally, we use a *class embedding* function  $\psi : \Omega \rightarrow \mathbb{R}^m$ , its inverse  $\psi^{-1}$ , and a function  $M : \mathbb{R}^n \rightarrow \mathbb{R}^m$  to map the two vector spaces so  $D(s) = \psi^{-1}(M(\xi(s)))$ .

$$S \xrightarrow{\xi} \mathbb{R}^n \xrightarrow{M} \mathbb{R}^m \xrightarrow{\psi^{-1}} \Omega \quad (3)$$

In this paper we explore three sentence embedding methods to implement  $\xi$ . We use a two-layer neural network as  $M$  and employ the standard *Mean Square Error* (MSE) as the inverse  $\psi^{-1}$ , to determine the closest embedding of each class  $\omega_i \in \Omega$  to the output of  $M$ .

### 4.2 Adapting Kartsaklis Method (LSTM)

Our basic inspiration for the algorithms of this paper is a text classification method proposed in (Kartsaklis et al., 2018) for the inverse dictionary problem where text definitions of terms are mapped to the term they define. The embedding of the class set into the continuous vector space (equivalent to the  $\psi$  function in equation 3) is done by expanding the knowledge graph of the dictionary words with nodes corresponding to words related to those terms. Next, random walks are performed on the graph to compute graph embeddings related to each dictionary node, using the *DeepWalk* algorithm (Perozzi et al., 2014).

A *Long Short-Term Memory* (LSTM) neural network, composed of two layers and an attention mechanism, is used in (Kartsaklis et al., 2018) for mapping the input texts to the input embedding vector space. To map the two continuous vector spaces representing the definitions and the dictionary terms, a two-layer neural network  $M$ , learned from the training dataset, is used.

For this work, the knowledge graph is replaced by an *intent proto-taxonomy*  $G$  which associates each class to a node and connects to them nodes which correspond to meta-knowledge concepts related to the class. To better capture the sequential



aspect of the intent proto-taxonomies, we also connect each class node to bigrams of concepts, i.e., the concatenation of two subsequent concepts. We represent this by the function  $\zeta$ , such as  $\zeta(\Omega) = G$ , which is invertible. Substituting this in equation 3,

$$S \xrightarrow{LSTM} \mathbb{R}^n \xrightarrow{M} \mathbb{R}^m \xrightarrow{DeepWalk} G \xrightarrow{\zeta^{-1}} \Omega \quad (4)$$

In practice, we compute the mapping from the class embedding space into the class set, called here  $InvG : \mathbb{R}^m \rightarrow \Omega$ , simply by determining the distance  $d$  between the output point in  $\mathbb{R}^m$  and the inverted projection of each class from  $\Omega$  and then considering the closest class. That is, for each  $w_i \in \Omega$ , we consider the associated node in  $G$  and compute the mapping in  $\mathbb{R}^m$  of that node:

$$InvG(x) = \arg \min_{w_i} d(x, DeepWalk(G(w_i))) \quad (5)$$

By substituting this function into equation 4, we obtain the algorithm we call here  $LSTM+T$ :

$$S \xrightarrow{LSTM} \mathbb{R}^n \xrightarrow{M} \mathbb{R}^m \xrightarrow{InvG} \Omega \quad (6)$$

For comparison, the traditional corresponding classification method is tested, where the graph embedding and associated functions are replaced by  $softmax+argmax$ . We call this  $LSTM$ :

$$S \xrightarrow{LSTM} \mathbb{R}^n \xrightarrow{M} \mathbb{R}^c \xrightarrow{softmax} \Omega \quad (7)$$

### 4.3 An Alternative to LSTM: USE

Recently, several new general-purpose language models that can be used for computing sentence embeddings have been proposed, among them the *Universal Sentence Encoder (USE)* (Cer et al., 2018). Such an approach consists of a *transformer* neural network (Vaswani et al., 2017), trained on varied sources of data, such as *Wikipedia*, web news, web question-answer pages and discussion forum. USE has achieved state-of-the-art results in various tasks, so we decided to try it in our experiments as an alternative to the LSTM for the embedding of input sentences.

In this work we employed the version 3 of the multilingual USE<sup>2</sup>. By replacing LSTM with USE in eq. 6 we obtain algorithm  $USE+T$ :

$$S \xrightarrow{USE} \mathbb{R}^n \xrightarrow{M} \mathbb{R}^m \xrightarrow{InvG} \Omega \quad (8)$$

<sup>2</sup><https://tfhub.dev/google/universal-sentence-encoder-multilingual/3>

Like in the previous case, we also compute the USE algorithm with traditional discrete softmax outputs for comparison, called here simply  $USE$ :

$$S \xrightarrow{USE} \mathbb{R}^n \xrightarrow{M} \mathbb{R}^c \xrightarrow{softmax} \Omega \quad (9)$$

### 4.4 Alternatives to DeepWalk

To explore variants of algorithms for embedding the classes and also approaches which do not need to be trained from scratch and allow on-the-fly handling of meta-knowledge, we tried replacing DeepWalk with two different methods.

The first one consists of applying USE sentence embeddings also for the class embeddings, such as in eq. 10. To simplify notation,  $emb$  represents either LSTM or USE embeddings for the input text.

$$S \xrightarrow{emb} \mathbb{R}^n \xrightarrow{M} \mathbb{R}^m \xrightarrow{USE^{-1}} G \xrightarrow{\zeta^{-1}} \Omega \quad (10)$$

This approach is similar to the way DeepWalk works but instead of training the graph embeddings from scratch, the class embeddings are represented by the mean sentence embedding computed from different random walks starting in the class node. We name such methods  $LSTM+S$  and  $USE+S$ , for  $emb$  substituted by LSTM and USE, respectively.

Additionally, we also evaluate the replacement of DeepWalk by the *Convolutional Deep Structured Semantic Model* (CDSSM) proposed in (Chen et al., 2016), yielding the following algorithm where  $emb$  can be either LSTM or USE embeddings.

$$S \xrightarrow{emb} \mathbb{R}^n \xrightarrow{M} \mathbb{R}^m \xrightarrow{CDSSM^{-1}} G \xrightarrow{\zeta^{-1}} \Omega \quad (11)$$

The CDSSM model consists of a three-layer convolutional neural network trained for creating embeddings of intent identifiers represented as sentences. In this work, we input to CDSSM the sequence of concepts listed in the nameId of each intent. We refer to these algorithms as  $LSTM+C$  and  $USE+C$ , for  $emb$  being substituted with LSTM and USE, respectively.

An intuitive way to understand those methods is to consider  $USE+T$  using a taxonomy as if its concepts had just abstract meanings: only their relations matter. In comparison,  $USE+S$  considers the meaning of the concepts besides their relations, while  $USE+C$  regards each nameId as a sentence, almost as if the developer had inputted a written description of the intent.

## 4.5 Out-of-Scope Sample Detection

In this paper we are interested both in the problems of: (1) deciding whether an user utterance is *in-scope* (IS) or *out-of-scope* (OOS) of the system; and (2) determining to which class an IS utterance belongs. For the former, a rejection mechanism based on a pre-defined threshold is used since it can be easily applied to all of the methods described previously without the need neither for any specific training procedure nor OOS training data.

In detail, suppose that for each class  $\omega_i \in \Omega$  there is a score denoted  $\phi_i \in Z$ , where  $|Z| = |\Omega|$ . Given that  $\max(Z)$  represents the highest score associated to a class and that a rejection threshold  $\theta$  has been defined on a validation set, samples can be classified as OOS whenever  $\max(Z) < \theta$ . If so, they are simply rejected, i.e., no classification output is produced for them. Otherwise, the sample is considered as in-scope and the classification is conducted normally.

The scores in  $Z$  are represented either by the *softmax* probability for the traditional softmax-based methods or by the similarity of sentence and intent embeddings for the proposed three approaches. For the latter, the similarity is computed by means of the dot product between the two embeddings.

## 5 Metrics, Datasets, and Experiments

In this section we present the experiments to evaluate the three algorithms described in the previous section, using each of the input embeddings LSTM and USE. We explore the impact on intent recognition both in terms of classifying correctly utterances (IS accuracy) and of finding which utterances are not covered by the intents (OOS accuracy).

### 5.1 Evaluation metrics

We employ a commonly-used metric for OOS detection, *equal error rate* (EER) (Tan et al., 2019), which corresponds to the classification error rate when the threshold  $\theta$  is set to a value where *false acceptance rate* (FAR) and *false rejection rate* (FRR) are the closest. These two metrics are defined as:

$$FAR = \frac{\text{number of accepted OOS samples}}{\text{total of OOS samples}} \quad (12)$$

$$FRR = \frac{\text{number of rejected IS samples}}{\text{total of IS samples}} \quad (13)$$

In addition, *in-scope error rate* (ISER) is considered to report IS performance, i.e. the error rate considering only IS samples when  $\theta$  is set to zero,

similar to the *class error rate* in (Tan et al., 2019). This metric is important to evaluate whether the proposed classification methods are able to keep up with the performance of the baselines in the main classification task.

### 5.2 The Larson and Telco Datasets

During the development and initial testing of the algorithms, we used two English datasets for in-depth experimentation. The first is the publicly-available *Larson* dataset (Larson et al., 2019); the second is a private real-world chatbot dataset used by a telecommunications provider for customer care, called here the *Telco* dataset. In the Larson dataset, we created an intent proto-taxonomy by hand, expanding the original identifiers of intents. The goal of the adjustments was to avoid spurious interference from taxonomy shortcomings or errors in the results. The complete list of the created taxonomic description of intents is listed in the appendix C to allow the reproduction of our results and further experimentation. In the Telco dataset, we created by hand the intent proto-taxonomy.

In the Larson dataset there is a total of 22,500 IS exemplars, evenly distributed across 150 classes, where 18,000 were used for training and 4,500 for testing. We conducted a simulation of OOS detection with the IS exemplars by doing 5 random samplings where we took out 30 intents and 3,600 training exemplars. We trained only with the remaining 120 intents and 14,400 exemplars. The test was then conducted using all the non-used 4,500 exemplars, where the 3,600 associated to the trained classes were considered the IS samples and the remaining 900 became OOS samples.

The Telco dataset contains 4,093 exemplars and 87 intents. From those, 3,069 exemplars were used for training and 1,024 for testing. The OOS scenario was simulated by extracting different random samplings where 5 intents were removed. Given the smaller size of this dataset compared to Larson, we conducted 20 samplings instead of 5.

For both sets we considered the following setup defined after preliminary evaluations. For the LSTM-based methods, the input sentence embedding size was set to 150 and output embeddings to 200. DeepWalk walk sizes were set to 20 for LSTM+T and USE+T. For both USE- and softmax-based methods we trained a two-layer neural network with 800 hidden neurons for 50 epochs.

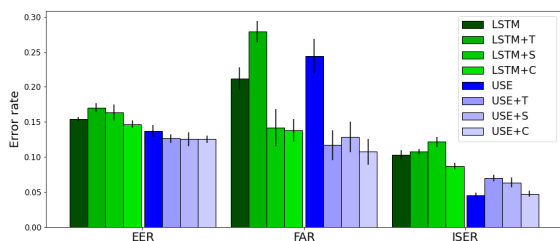


Figure 3: Different methods to incorporate the intent proto-taxonomy in Larson dataset, compared to the LSTM and USE baselines.

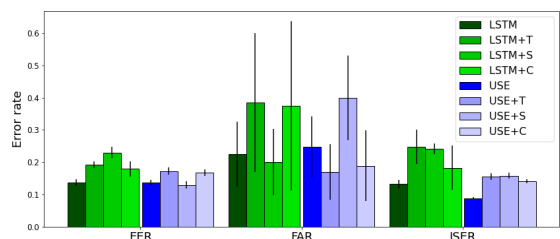


Figure 4: Different methods to incorporate the intent proto-taxonomy in Telco dataset, compared to the LSTM and USE baselines.

### 5.3 Results in the Larson and Telco Datasets

The results on the Larson dataset are graphically depicted in fig. 3. We observed that there was a slight improvement (a decrease) in EER, especially with the USE-based and the LSTM+C methods. More notably, there was a significant improvement in terms of FAR for all USE-based methods and LSTM+S and LSTM+C. Notice that even though the proposed approaches generally did not outperform LSTM and USE in ISER (except LSTM+C), we observed that the methods with better ISER tended to produce also better EER and FAR.

In fig. 4, we see that the results on the Telco dataset presented a different scenario. The proposed methods generally performed worse than or, at best, similar to LSTM and USE in EER. In terms of FAR, some methods such as USE+T and USE+C seem to outperform the others but, considering the high standard deviations, the results were not significant. On the other hand, we also observed that the methods failed to get close in ISER compared to the softmax-based methods. That seems to indicate that for the cases where making use of meta-knowledge harms too much ISER, the symbolic knowledge did not decrease neither EER nor FAR.

There were two key findings from our experiments with the Larson and the Telco datasets. First, the improvements using LSTM or USE as a base-

line seemed to be similar, possibly slightly better for the USE algorithm. Second, and most importantly, we saw much more improvement in the use of the intent proto-taxonomy in the Larson than in the Telco dataset, in spite of the similar nature of the datasets and the intent proto-taxonomies. This motivated us to try out the ideas in a larger and more diverse number of workspaces and solely focusing on USE to simplify the experiments.

### 5.4 The ChatWorks Dataset

To test our algorithms in a context of high diversity and realism, we used the same large set of real, professional workspaces explored in (Pinhanez et al., 2021), which come from the professional chatbot development platform *ChatWorks*.

We started with the 3,840 workspaces available in English. To eliminate possible problems due to workspaces with poor quality, we employed the  $3\sigma$ -rule, where values smaller greater than 3 standard deviations from the mean are not considered. Workspaces with the number of intents or exemplars below and above those thresholds were removed. Also, to avoid workspaces with few exemplars per intent, the ratio of the number of exemplars to the number of intents had to be greater than 10. From the filtered set we randomly selected 200 workspaces for testing.

The evaluation involved the execution of 20 iterations for each workspace. The tests were performed for all USE-based methods (USE, USE+T, USE+S, and USE+C). First, the workspaces were split into training and test datasets (75% and 25%, respectively). Next, the four methods were trained and tested on these datasets. The evaluation metrics (EER, FAR, and ISER) were then measured on the results for the test datasets and the average errors and their standard deviations were computed.

### 5.5 Results in the ChatWorks Dataset

Appendix D contains a table with the results for each of the 200 workspaces in the ChatWorks dataset. Figure 5 summarizes the results of the experiments showing the distribution of the 200 workspaces according to ranges of the improvement of each of the three methods compared to the baseline of USE. Improvement is calculated by subtracting the errors in each of our proposed methods from the errors in the USE baseline (error values are scores between 0 and 1). When one of our methods was worse than the baseline then  $diff < 0$ ,

since smaller is better, and conversely for when it is better than the baseline, i.e.,  $diff \geq 0$ .

The results shown in fig. 5 indicate that the USE+C algorithm achieved the best results in all three metrics, although there is a significant portion of workspaces where the other methods also did well, especially in OOS detection (FAR).

But, more important, the results seem to support our claim that meta-knowledge embedded in the output layer of our neuro-symbolic algorithms can improve intent recognition performance in practical systems. Notably in OOS detection (FAR), 67% of the workspaces experienced a decrease in the error rate using USE+C. Besides, in 39% of the workspaces we observed a decrease in the error rate of more than 0.05 (in a 0 to 1 scale), and in 23%, of more than 0.1. The USE+T also did well with similar but slightly smaller decreases in error.

Overall, the error rates for the EER metric also decreased in relation to the baseline. Figure 5 shows that 41% of the workspaces had some level of decrease in EER with the USE+C algorithm, in 10% of them with decreases of 0.05 or more. However, the results for the in-scope accuracy (ISER) were much smaller with only about 16% of the workspaces having any kind of decrease.

The ChatWorks dataset, as noted before, includes all kinds of workspaces. Taxonomy rates varies anywhere from 0 to 1, and there are very small and very large workspaces. To test our methods in a scenario closer to a professional, well-developed chatbot, we filtered further the dataset to include only workspaces with taxonomy rate greater or equal to 0.7, with number of intents equal or more than 32, and at least an average of 25 exemplars per intent, resulting in 18 workspaces.

Figure 6 shows the distribution of the results of the experiments with those 18 workspaces, which were better than in the full ChatWorks dataset. Both USE+T and USE+C yielded EER decreases in 50% or more of the workspaces. Moreover, 83% of the workspaces decreased the FAR error, either with USE+T or USE+C, and both decreased FAR in more than a third of the workspaces by more than 0.1. We discuss the results and implications next.

## 6 Discussion and Future Work

We started this paper by proposing the combination of exemplars and symbolic characterizations of a class as a way to enhance ML-based intent recognition. We proposed 3 new neuro-symbolic algo-

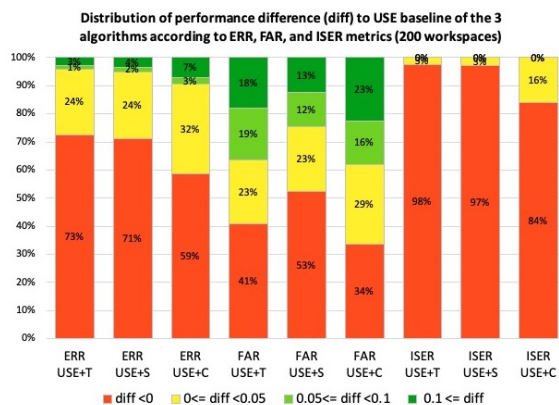


Figure 5: Distribution of performance difference (diff) to USE baseline of the 3 methods according to EER, FAR, and ISER metrics in all 200 workspaces of the Chatworks dataset.

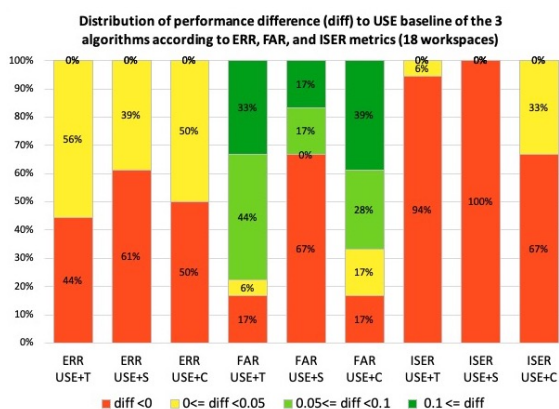


Figure 6: Distribution of performance difference (diff) to USE baseline of the 3 methods according to EER, FAR, and ISER metrics in the 18 most developed workspaces of the Chatworks dataset.

gorithms and tested them using datasets built using data from intent identifiers of conversational systems. Such identifiers often store taxonomic-like structures, due to a common practice among developers of professional conversational systems (Pinhanez et al., 2021). The results of the experiments indicate that the intent proto-taxonomies embedded by those developers can indeed be used by many workspaces to improve accuracy in intent recognition, notably in OOS detection.

We see as one of the main contributions of this paper the creation of methods with which ML engineers can improve the accuracy of their systems by simply mining “documentation” from chatbots, without any further data and annotation.

Our results show that almost 40% of the 200 professional workspaces drawn from ChatWorks saw decreases of more than 0.05 in OOS detection



error rates. Also, in 42% of them the overall error rate was improved, using the USE+C algorithm. When considering the more well-structured and developed 18 workspaces, we saw much higher gains with the USE+T algorithm. Those accuracy improvements were achieved without any change in the training set but simply by incorporating the meta-knowledge into intent recognition.

Notice that the testing methodology used in this work is considerably harder than the practice of the majority of research papers, since it evaluates performance in 200 professional, non-edited workspaces from different domains. In reality, most ML algorithms do not perform well in all datasets, and ML practitioners often test different algorithms and parameters until accuracy is good enough.

However, the improvement in OOS detection (FAR) was not mirrored in classification error (ISER). First, we must keep in mind that intent classification is often performed in two steps, first OOS sentences detection and removal, followed by intent classification of the IS sentences. Given the improvements observed in OOS detection, it would make sense to use our algorithms in the first step for many of the ChatWorks workspaces (about 60% of them), and selectively use it for IS classification only when it works better than the baseline.

But why were there so many workspaces where we did not see impact? It is important to take into account that the ChatWorks dataset has workspaces in different stages of development and deployment. By selecting better quality workspaces, we saw much higher gains. We explored briefly characterizations of the intent proto-taxonomy quality, such as taxonomy rate, depth of the taxonomy, and number of concepts, but we saw no clear correlation with decreases in error rates. We believe more complex metrics of knowledge structure need to be employed to characterize which intent proto-taxonomies are likely to have the greatest impacts. We plan to do so in our future work.

It is important to notice that, in the workspaces where we did see impact, the symbolic knowledge was mined from an absolutely “raw” format. In spite of that, by using the basic graph mining method described in appendix A, it was possible to obtain a “meaningful” taxonomic structure, similar to a knowledge graph which could be used by our neuro-symbolic algorithms. To improve the quality of the taxonomies, we are working on designing an interface which allows the developers to manipu-

late directly the intent proto-taxonomy to make it more correct and complete, so to possibly decrease even more the intent recognition error rates.

We have demonstrated in this work that combining exemplar and symbolic ways of defining classes can have a positive impact in the performance of the recognition system. This was done in the context of conversational systems where developers fortuitously embed such alternative descriptions of classes in their name identifiers. We believe it is possible to find in other machine learning development platforms similar patterns of knowledge embedding.

For example, we know that it is common for people to use similar taxonomic structures when naming file and e-mail folders, giving names to functions and variables in programs and data, and writing comments into *Jupyter* notebooks. Also, ML platforms can further foster the use of metadata by developers by explicitly asking them to input, besides exemplars, categorical or textual descriptions of the classes.

As we move along the path of creating such neuro-symbolic systems, not only we should expect that the job of developers becomes easier, as they follow their own cultural and linguistic practices, but also that machines became better in recognizing those classes accurately. Using multiple forms of class definitions can be a winning proposition for both ML systems and their developers.

## Ethical Issues

The ChatWorks dataset was composed only of workspaces in which the developers explicitly opted-in to share their code and content for research and development purposes with the company which owns the platform. Those workspaces were shared by the company with the authors of this paper with a clear condition of not publicly sharing their contents and publishing only aggregated results or in an anonymous form. We do not see any specific impact of those limitations in the results of our research but they preclude easy forms of replication of our results with that dataset. To better enable reproducibility, we presented the analysis of the public Larson dataset and shared the intent proto-taxonomy we created manually from its original intent structure in appendix C.

## References

- Zeynep Akata, Florent Perronnin, Zaid Harchaoui, and Cordelia Schmid. 2016. [Label-embedding for image classification](#). *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(7):1425–1438.
- Zeynep Akata, Scott Reed, Daniel Walter, Honglak Lee, and Bernt Schiele. 2015. Evaluation of output embeddings for fine-grained image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yoshua Bengio. 2017. [The consciousness prior](#).
- Tarek R. Besold, Artur d’Avila Garcez, Sebastian Bader, Howard Bowman, Pedro Domingos, Pascal Hitzler, Kai-Uwe Kuehnberger, Luis C. Lamb, Daniel Lowd, Priscila Machado Vieira Lima, Leo de Penning, Gadi Pinkas, Hoifung Poon, and Gerson Zaverucha. 2017. [Neural-symbolic learning and reasoning: A survey and interpretation](#).
- Geoffrey C Bowker and Susan Leigh Star. 2000. *Sorting things out: Classification and its consequences*. MIT press.
- Paulo Cavalin, Victor Henrique Alves Ribeiro, Ana Appel, and Claudio Pinhanes. 2020. [Improving out-of-scope detection in intent classification by using embeddings of the word graph space of the classes](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3952–3961, Online. Association for Computational Linguistics.
- Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Brian Strope, and Ray Kurzweil. 2018. [Universal sentence encoder for English](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 169–174, Brussels, Belgium. Association for Computational Linguistics.
- Y. Chen, D. Hakkani-Tür, and X. He. 2016. Zero-shot learning of intent embeddings for expansion by convolutional deep structured semantic models. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6045–6049.
- Andrea Civan, William Jones, Predrag Klasnja, and Harry Bruce. 2008. [Better to organize personal information by folders or by tags?: The devil is in the details](#). *Proceedings of the American Society for Information Science and Technology*, 45(1):1–13.
- Luc De Raedt, R. Manhaeve, S. Dumancic, Thomas De-meester, and A. Kimmig. 2019. [Neuro-symbolic = neural + logical + probabilistic](#). In *Proceedings of the 2019 International Workshop on Neural-Symbolic Learning and Reasoning*, page 4.
- Emile Durkheim and Marcel Mauss. 1963. *Primitive Classification*, volume 273. University of Chicago Press.
- Lisa Ehrlinger and Wolfram Wöß. 2016. Towards a definition of knowledge graphs. *SEMANTiCS (Posters, Demos, SuCCESS)*, 48:1–4.
- Marco Fossati, Dimitris Kontokostas, and Jens Lehmann. 2015. Unsupervised learning of an extensive and usable taxonomy for dbpedia. In *Proceedings of the 11th International Conference on Semantic Systems, SEM ’15*. ACM.
- Artur d’Avila Garcez, Marco Gori, Luis C. Lamb, Luciano Serafini, Michael Spranger, and Son N. Tran. 2019. Neural-symbolic computing: An effective methodology for principled integration of machine learning and reasoning. *Journal of Applied Logics*, 6(4):611–631.
- Drew Hudson and Christopher D Manning. 2019a. Learning by abstraction: The neural state machine. In *Advances in Neural Information Processing Systems*, pages 5903–5916.
- Drew Hudson and Christopher D Manning. 2019b. [Learning by abstraction: The neural state machine](#). In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 5903–5916. Curran Associates, Inc.
- Shaoxiong Ji, Shirui Pan, Erik Cambria, Pekka Martinen, and Philip S. Yu. 2020. [A survey on knowledge graphs: Representation, acquisition and applications](#).
- Daniel Jurafsky and James H. Martin. 2009. *Speech and Language Processing (2nd Edition)*. Prentice-Hall, Inc., USA.
- Dimitri Kartsaklis, Mohammad Taher Pilehvar, and Nigel Collier. 2018. [Mapping text to knowledge graph entities using multi-sense LSTMs](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1959–1970, Brussels, Belgium. Association for Computational Linguistics.
- Stephen M Kosslyn. 2006. You can play 20 questions with nature and win: Categorical versus coordinate spatial relations as a case study. *Neuropsychologia*, 44(9):1519–1523.
- Stefan Larson, Anish Mahendran, Joseph J. Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K. Kummerfeld, Kevin Leach, Michael A. Laurenzano, Lingjia Tang, and Jason Mars. 2019. [An evaluation dataset for intent classification and out-of-scope prediction](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1311–1316, Hong Kong, China. Association for Computational Linguistics.

- Robin Manhaeve, Sebastijan Dumancic, Angelika Kimmig, Thomas Demeester, and Luc De Raedt. 2018. [Deepproblog: Neural probabilistic logic programming](#). In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 3749–3759. Curran Associates, Inc.
- Christopher Manning and Hinrich Schütze. 1999. *Foundations of statistical natural language processing*. MIT press.
- Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B. Tenenbaum, and Jiajun Wu. 2019. [The Neuro-Symbolic Concept Learner: Interpreting Scenes, Words, and Sentences From Natural Supervision](#). In *International Conference on Learning Representations*.
- Rodney Needham. 1979. *Symbolic classification*. Goodyear Publishing Company.
- Allen Newell. 1973. You can't play 20 questions with nature and win: Projective comments on the papers of this symposium. In W.G. Chase, editor, *Visual information processing*. Academic Press.
- Alessandro Oltramari, Jonathan Francis, Cory Henson, Kaixin Ma, and Ruwan Wickramarachchi. 2020. [Neuro-symbolic architectures for context understanding](#).
- Mark Palatucci, Dean Pomerleau, Geoffrey E Hinton, and Tom M Mitchell. 2009. [Zero-shot learning with semantic output codes](#). In *Advances in Neural Information Processing Systems*, volume 22. Curran Associates, Inc.
- Emilio Parisotto, Abdel-Rahman Mohamed, Rishabh Singh, Lihong Li, Dengyong Zhou, and Pushmeet Kohli. 2017. Neuro-symbolic program synthesis. In *International Conference on Learning Representations (ICLR)*.
- Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 701–710.
- Claudio Santos Pinhanez, Heloisa Candello, Paulo Cavalin, Mauro Carlos Pichiliani, Ana Paula Appel, Victor Henrique Alves Ribeiro, Julio Nogima, Maira de Bayser, Melina Guerra, Henrique Ferreira, et al. 2021. Integrating machine learning data with symbolic knowledge from collaboration practices of curators to improve conversational systems. In *Proceedings of the 2021 ACM Conference on Human Factors in Computing Systems (CHI'21)*, pages 1–13.
- Whitman Richards. 1982. How to play twenty questions with nature and win.
- Richard Socher, Milind Ganjoo, Christopher D Manning, and Andrew Ng. 2013. [Zero-shot learning through cross-modal transfer](#). In *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.
- Ming Tan, Yang Yu, Haoyu Wang, Dakuo Wang, Saloni Potdar, Shiyu Chang, and Mo Yu. 2019. [Out-of-domain detection for low-resource text classification tasks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3566–3572, Hong Kong, China. Association for Computational Linguistics.
- Joshua B. Tenenbaum, Charles Kemp, Thomas L. Griffiths, and Noah D. Goodman. 2011. [How to grow a mind: Statistics, structure, and abstraction](#). *Science*, 331(6022):1279–1285.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Wei Wang, Vincent W Zheng, Han Yu, and Chunyan Miao. 2019. A survey of zero-shot learning: Settings, methods, and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2):1–37.
- Steve Whittaker, Tara Matthews, Julian Cerruti, Hernan Badenes, and John Tang. 2011. [Am I wasting my time organizing email? a study of email refinding](#). In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI '11*, page 3449–3458, New York, NY, USA. Association for Computing Machinery.
- Bai Yang, Zhang Liping, and Zhao Fengrong. 2019. [A survey on research of code comment](#). In *Proceedings of the 2019 3rd International Conference on Management Engineering, Software Engineering and Service Sciences, ICMSS 2019*, page 45–51, New York, NY, USA. Association for Computing Machinery.

## A An Algorithm to Extract Intent Proto-Taxonomies from NameIds

In our previous work (Pinhanez et al., 2021) we described an algorithm to mine proto-taxonomies from the nameIds of a workspace. For completeness, we include it here. It consists of three steps:

1. Finding the best separator to split the name into a sequence of *concepts*;
2. Splitting the nameIds with the selected separator.
3. Generating an intent proto-taxonomy using the terms split by the best separator as concepts and considering consecutive concepts in a nameId as having a link between them.

In order to find the best separator, our algorithm first calculates the *perplexity* of the bag of concepts using each separator. Perplexity is a measure of uncertainty for a given sequence of words (or concepts) appearing in a language model (Manning and Schütze, 1999). For that, we build probabilistic language models based on *bigrams* using the *maximum likelihood estimator* as in (Jurafsky and Martin, 2009, Chapter 3), and then compute the average perplexity using a standard *leave-one-intent-out* evaluation scheme. The separator which minimizes perplexity is chosen as the separator for the workspace.

## B Intent Proto-Taxonomies as a Common Practice in ChatWorks

In (Pinhanez et al., 2021) we evaluated the use of proto-taxonomies by the curators of the ChatWorks dataset, by developing a metric, called *taxonomy rate*, which is the ratio of the number of nameIds which have a proto-taxonomy embedded in it to the total number of intents of the workspace.

To determine whether a given nameId has a proto-taxonomy embedded in it, we considered three criteria: (1) the nameId has two or more concepts; (2) there is at least one other nameId which has at least one identical concept at the exact same level; and (3) the concept does not appear in all nameIds of the workspace in that level. After all the nameIds of a workspace were determined as having an embedded intent or not, the *taxonomy rate* was calculated by considering the ratio between the number of intents with a taxonomic structure and the number of intents of the workspace.

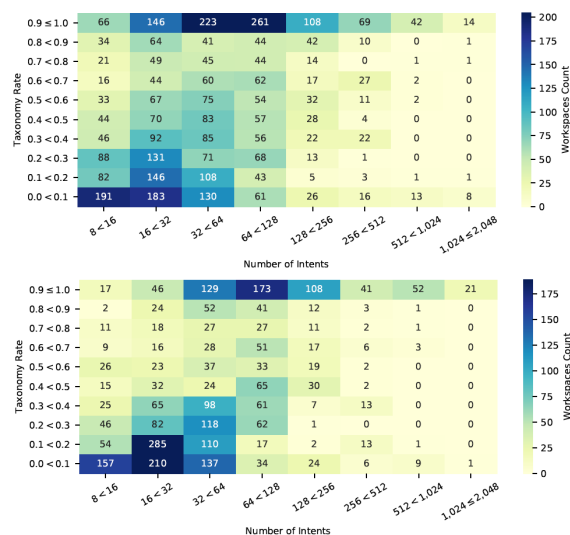


Figure 7: Distribution of the number of workspaces according to the taxonomy rate of the 3,840 English (top) and 2,895 Portuguese (bottom) workspaces. Notice that the X axis is in logarithmic scale.

A taxonomy rate of above 10% was observed in 76% of the English workspaces and in 80% the Portuguese ones. Almost 52% of all workspaces, in English, and 38% in Portuguese, had a taxonomy rate larger than 50%. Moreover, 16% and 20% of the workspaces, for English and Portuguese respectively, had a very high taxonomy rate from 90% to 100%. Considering both the English and Portuguese workspaces, about 70% of them have a taxonomy rate above 20% and more than half of the workspaces above 64 intents have a taxonomy of 50% or more.

Figure 7 shows how the workspaces are distributed considering both the number of intents (X axis) and the taxonomy rate (Y axis). Notice that the distribution, in both languages, follows a sort of a “step” function where, as the threshold between 32 and 64 intents is crossed, the majority of the workspaces has more than 50% of taxonomy rate. More details can be found in (Pinhanez et al., 2021).

## C The Intent Proto-Taxonomy Created for the Larson Dataset

In tables 1 and 2, we list the taxonomy which was manually created for the Larson dataset, with the original nameId on the left and the created taxonomic representation on the right, represented as a string of concepts separated by spaces.



nameId	Concepts
accept_reservations	accept reservation
account_blocked	account blocked
alarm	set alarm
application_status	application status
apr	what month
are_you_a_bot	you bot
balance	what balance
bill_balance	what bill balance
bill_due	when bill due
book_flight	book flight
book_hotel	book hotel
calculator	calculate
calendar	calendar
calendar_update	calendar update
calories	calories dish
cancel	cancel action
cancel_reservation	cancel reservation
car_rental	car rental
card_declined	card declined
carry_on	carry-on rule
change_accent	change accent
change_ai_name	change bot name
change_language	change language
change_speed	change speed
change_user_name	change user name
change_volume	change volume
confirm_reservation	confirm reservation
cook_time	cook time
credit_limit	credit limit
credit_limit_change	credit limit change
credit_score	credit score
current_location	what current location
damaged_card	damaged card
date	what date
definition	definition
direct_deposit	direct deposit
directions	what direction
distance	what distance
do_you_have_pets	do you have pet
exchange_rate	exchange rate
expiration_date	expiration date
find_phone	find phone
flight_status	flight status
flip_coin	flip coin
food_last	food last
freeze_account	block account
fun_fact	fun fact
gas	gas level
gas_type	gas type
goodbye	goodbye
greeting	greeting
how_busy	how busy
how_old_are_you	how old you
improve_credit_score	improve credit score
income	what income
ingredient_substitution	ingredient substitution
ingredients_list	ingredient list
insurance	insurance benefit
insurance_change	insurance change
interest_rate	interest rate
international_fees	international fee
international_visa	international visa
jump_start	jump start
last_maintenance	last maintenance
lost_luggage	lost luggage
make_call	make call
maybe	maybe
meal_suggestion	meal suggestion
meaning_of_life	what meaning life
measurement_conversion	measurement conversion
meeting_schedule	meeting schedule
min_payment	minimum payment
mpg	what mpg
new_card	apply card
next_holiday	next holiday
next_song	next song
no	no
nutrition_info	nutrition info
oil_change_how	how change oil
oil_change_when	when change oil
order	order shopping
order_checks	order check
order_status	order status
pay_bill	pay bill
payday	when payday
pin_change	change pin
play_music	play music
plug_type	what plug type
pto_balance	pto balance
pto_request	pto request
	...

Table 1: The taxonomy created for the Larson dataset.

nameId	Concepts
	...
pto_request_status	pto request status
pto_used	pto used
recipe	recipe dish
redeem_rewards	redeem reward
reminder	reminder action
reminder_update	reminder update
repeat	repeat action
replacement_card_duration	replacement card duration
report_fraud	report fraud
report_lost_card	report lost card
reset_settings	reset setting
restaurant_reservation	restaurant reservation
restaurant_reviews	restaurant review
restaurant_suggestion	restaurant suggestion
rewards_balance	reward balance
roll_dice	roll dice
rollover_401k	rollover 401k
routing	find routing
schedule_maintenance	schedule maintenance
schedule_meeting	schedule meeting
share_location	share location
shopping_list	shopping list
shopping_list_update	shopping list update
smart_home	smart home
spelling	spelling word
spending_history	spending history
sync_device	sync device
taxes	what taxes
tell_joke	tell joke
text	text person
thank_you	thank
time	what time
timer	set timer
timezone	set timezone
tire_change	tire change
tire_pressure	tire pressure
todo_list	todo list
todo_list_update	todo list update
traffic	what traffic
transactions	card transaction
transfer	transfer account
translate	translate word
travel_alert	travel alert
travel_notification	travel notification
travel_suggestion	travel suggestion
uber	get uber
update_playlist	update playlist
user_name	user name
vaccines	what vaccine
w2	get w2
weather	what weather
what_are_your_hobbies	what hobby
what_can_i_ask_you	what ask you
what_is_your_name	what name
what_song	what song
where_are_you_from	where you from
whisper_mode	whisper mode
who_do_you_work_for	who you work
who_made_you	who made you
yes	yes
	...

Table 2: The taxonomy created for the Larson dataset (cont.).

## D Results of the Proposed Algorithms in the ChatWorks Dataset

Table 3 shows the results of the proposed three algorithms and the baseline for each workspace in the ChatWorks dataset. For each metrics, EER, FAR, and ISER, the table shows the average accuracy and the standard deviation over the experiments made with the 20 random splits. The table also lists information characterizing the workspaces: number of intents, number of exemplars, the average number of exemplars per intent, taxonomy rate, the average depth of the taxonomy graph, and the number of different concepts.

